# GHAJAR
# EXHIBIT 30

Page 1

UNITED STATES DISTRICT COURT

NORTHERN DISTRICT OF CALIFORNIA

SAN FRANCISCO DIVISION

_____

RICHARD KADREY, et al.,            )
                                   )
            Individual and         )
            Representative         )
            Plaintiffs,            )
                                   )
v.                                 )    Case No. 3:23-cv-03417-VC
                                   )
META PLATFORMS, INC.,              )
                                   )
            Defendant.             )
_____)


   ** HIGHLY CONFIDENTIAL - ATTORNEYS' EYES ONLY **

         Videotaped Deposition of SEAN BELL

            San Francisco, California

           Wednesday, December 11, 2024



           Reported Stenographically by

    Michael P. Hensley, RDR, CSR No. 14114


_____

              DIGITAL EVIDENCE GROUP

           1730 M. Street, NW, Suite 812

              Washington, D.C. 20036

                (202) 232-0646

Page 12

```
 1        Q.   All right.   And what's your position

 2   there?

 3        A.   I'm a research scientist manager.

 4        Q.   All right.   And when did you start at Meta

 5   Platforms?

 6        A.   In January 2019.

 7        Q.   And when you started at Meta Platforms in

 8   2019, did you hold the same position that you hold

 9   today?

10        A.   Yes.   I -- my job offer may have said

11   something different.   I don't remember the exact

12   title, but I've -- I've been a manager since I

13   started.

14        Q.   Okay.

15        A.   It may have said "research scientist," for

16   example.

17        Q.   Have you been promoted since you began

18   working at Meta?

19        A.   Yes.   I was promoted from Level 6 to

20   Level 7.

21        Q.   And what does that mean to be promoted

22   from Level 6 to Level 7?
```

12/11/2024          Richard Kadrey, et al. v. Meta Platforms, Inc.          Sean Bell
Highly Confidential - Attorneys' Eyes Only

Page 24

```
 1       A.   Yes.  That's one of the things that I do.

 2       Q.   Okay.  And you said that you work on --

 3  you also work on collecting data for other

 4  foundation models.

 5            Did I hear that correctly?

 6       A.   That's right.

 7       Q.   What other foundation models do you

 8  collect data for?

 9       A.   So I'm currently collecting data for

10  Llama 4 for the -- so that's the first one.  The

11  second category is the image generation models.  We

12  called it "Emu," E-m-u.  And then the third category

13  is the movie generation models called "Movie Gen."

14            And so sometimes the data that my team

15  collects, even though it was intended for --

16  primarily for one of those three categories of

17  models, it might be useful to other teams; and so

18  occasionally we'll establish collaborations where we

19  might additionally share data with other teams.

20       Q.   Okay.  And in your role working on the

21  Llama models, do you provide any strategic guidance

22  on collecting or using data?
```

Page 25

1          A.    Yes.

2                ATTORNEY HARTNETT:    Object to the form.

3                THE WITNESS:   Yes.

4                          ///

5    BY ATTORNEY STOLER:

6          Q.    Okay.   What sort of guidance do you

7    provide as to the use of data for training Llama?

8          A.    So ultimately what my team is accountable

9    for is the final dataset that goes into Llama.   And

10   what -- what matters the most is, you know, what is

11   the benchmark performance as a result of training on

12   this data?

13               And so our goal is to maximize the

14   benchmark performance as a result of training on the

15   data, and so there's a huge number of different

16   steps prior to that point, prior to the actual data

17   you put into Llama.

18               And so I might be consulting on any part

19   of the strategy that leads up to that point.   For

20   example, how do you process data?   How do you curate

21   data?   How do you measure data quality?   How can you

22   decide that one dataset is better than another

Page 26

1    dataset?  You know, what are the -- what is even the

2    scientific method to know how to attribute, you

3    know, this concept of, quote/unquote, "quality"?

4    Which means, really, the benchmark performance as a

5    result of training on the data.  You know, it may

6    also be what -- you know, thinking about what kinds

7    of sources we need to be collecting.

8            And so generally -- either I might be

9    writing some of the strategies, or people -- or

10   leadership on my team may be writing those

11   strategies and then I'm reviewing it.

12   BY ATTORNEY STOLER:

13       Q.   And prior to coming to Meta, did you ever

14   work on large language models?

15       A.   No.

16       Q.   So fair to say you learned about data

17   strategy for LLMs while at Meta?

18       A.   Yes.

19       Q.   Okay.  And who did you learn that from?

20       A.   For LLMs specifically, I read a lot of

21   academic papers.  And also I -- you know, as a

22   manager, I took on a number of senior individual

12/11/2024          Richard Kadrey, et al. v. Meta Platforms, Inc.          Sean Bell
Highly Confidential - Attorneys' Eyes Only

Page 41

```
 1   BY ATTORNEY STOLER:

 2        Q.   Let me rephrase.

 3        A.   Yeah.

 4        Q.   Does the quality of the pretraining data

 5   affect the performance of an LLM?

 6             ATTORNEY HARTNETT:   Object to the form.

 7             THE WITNESS:   So what I was trying to get

 8   at earlier is that we've basically defined our

 9   notion of quality to be the performance of the LLM,

10   and so the benchmarks are a way to estimate the --

11   you know, estimate the performance.

12             You know, this was actually a major

13   challenge and question that we had to think about

14   when we started the team, which is how do you think

15   about quality and is it possible to know what it is.

16   And so we had some heuristics but we didn't have

17   good science connecting it.

18             And so a lot of the work that this team is

19   doing is trying take something that's a little bit

20   subjective and turn it into something as scientific

21   as possible, meaning that, you know, there's no sort

22   of judgment around it.   It's a connection between a
```

12/11/2024          Richard Kadrey, et al. v. Meta Platforms, Inc.          Sean Bell
Highly Confidential - Attorneys' Eyes Only

Page 42

```
 1    particular experiment that you do and then a very

 2    measurable set of numbers.

 3              (Clarification inquiry by The Court

 4              Reporter.)

 5                          ///

 6    BY ATTORNEY STOLER:

 7         Q.    What does the term -- in the context of

 8    pretraining LLMs, what does the term "data mix"

 9    refer to?

10         A.    So "data mix" refers to the final overall

11    combination of data that the LLM is actually trained

12    on.  And so there's really two kinds of activities

13    when creating -- when collect -- when creating a

14    dataset.

15              One of them will be on an individual

16    dataset basis, which types or sources of data do we

17    want to use, and which subset of the data -- or how

18    do I filter it or process it?  And so that may give

19    us, say, order of approximately 100 different

20    individual sources.

21              The data mix is an additional step where I

22    now need to choose the relative weighting or
```

Page 43

1    proportion of those individual datasets.  And so if

2    the weight is particularly high, I might repeat

3    certain datasets more often.  If the weight is less

4    than 1, for example, then I might be taking a

5    proportional subsets; maybe I'll sample only 80

6    percent of some data source, for example.

7              And so the data mix is the work of

8    deciding all of those weights and proportions and

9    maybe whether or not to include a dataset at all.

10        Q.    Do you oversee the determination of the

11   data mix for Llama 4's pretraining data?

12        A.    Yes.

13        Q.    Help me understand.

14              You testified that you are not aware of

15   what data goes into -- sorry, let me -- I'll start

16   again.

17              Help me understand.  Did you previously

18   mention that you don't know what datasets were used

19   to train Llama 3 but are not being used to train

20   Llama 4?

21              ATTORNEY HARTNETT:  Objection to form.

22   Asked and answered.

Page 49

1    Llama 4 pretraining dataset, would that render your

2    picture of the training data incomplete?

3                   ATTORNEY HARTNETT:   Objection to the form.

4                   THE WITNESS:   I'm not sure what you mean

5    by that.

6    BY ATTORNEY STOLER:

7         Q.    Yeah.   If you don't know what datasets

8    went into pretraining Llama 4, would you be able to

9    evaluate whether the dataset you created was

10   optimal?

11                  ATTORNEY HARTNETT:   Objection to the form.

12                  THE WITNESS:   So we -- so a lot of what

13   we're focused on -- I mean, I'm -- honestly, I'm not

14   sure that it's needed.

15                  So I mean, we -- what we do is we -- we've

16   spent a lot of time figuring out what is the science

17   of what makes for good data in terms of the

18   benchmark performance.   And so we have done a lot of

19   innovations in this space and come up with new types

20   of metrics and ways to understand the data.

21                  And some of what we've learned in the past

22   is that human intuition on data is -- is wrong.   And

Page 50

```
 1   so in 2023, when I was not there, but I -- you know,

 2   other researchers --

 3              I'm sorry.  My microphone dropped.  Let me

 4   put it back, sorry.

 5              In 2023, I was told about research

 6   projects where people attempted to look at the data

 7   and make human judgements just on looking at the

 8   data, did they think it would be useful to LLM

 9   training or not.

10              And in a different workstream, they said

11   let's take a statistical approach where we only look

12   at the benchmark performance, and what they found

13   was that the human intuition actually led you to a

14   worse-performing LLM in the end and that it was

15   counterproductive to attempt to, on a per-document

16   or per-dataset basis, make human judgments about

17   that.

18              And so when I say that I've "shifted," it

19   was based on scientific learning; that the best way

20   to do data mix optimization is to, as much as

21   possible, distill everything into a statistical

22   view.
```

Page 51

1                Again, when I say "statistical view," I

2       mean, like, the distribution of tokens, the

3       distribution of benchmark scores, and, you know,

4       look at that view as much as possible, make all your

5       decisions in that lens.   And then separately have a

6       central private -- you know, a privacy

7       infrastructure that guarantees that we're following

8       all of the -- all of our mitigation decisions and

9       that -- you know, one of the principles I've -- I've

10      pushed the team building that system was this idea

11      that you couldn't train on bad data even if you

12      tried.   And when I say "bad data," I mean data that

13      disagrees with any of our privacy decisions.

14      BY ATTORNEY STOLER:

15          Q.   You mention human -- let me back up.

16               You just said that there was a switch in

17      the approach that you took to evaluating data.

18               Did I get that correct?

19          A.   I wouldn't call it a switch.   And so,

20      again, I -- I wasn't there; so I've only heard, you

21      know, small pieces of what happened.

22          Q.   Hmm.

Page 52

1     A.     But my understanding, hearing from others,

2   is that there is a team in Europe that questioned

3   whether or not the team in the U.S. was doing the

4   right thing in terms of pretraining data.  And so

5   they started a separate project to do human

6   annotation of data to try and use human judgments to

7   label each dataset -- is this a good quality dataset

8   or not? -- and, you know, build a whole process

9   around this.

10          You know, basically what happened was the

11  team in Europe looked at the pretraining data

12  directly, and they -- they thought, "This is very

13  low-quality data.  You guys don't know what you're

14  doing.  How -- why are we training on such bad data

15  from this human intuition?"

16          And so they started a project to do

17  annotation of this data to then label data as "good

18  quality" or "bad quality."  And then it turns out

19  that that whole approach of annotation did not lead

20  to more useful filters, and so we're not using that

21  approach.

22          So it wasn't so much that we switched

Page 53

1    approaches.   It was that there were a lot of

2    different approaches being tried at the same time,

3    and one of them was this annotation process to

4    challenge the results of the other team.

5              You know, and the last thing I'll add is

6    that I wasn't here for any of this, and so this is

7    me, you know, recalling what -- how other people

8    described to me how it went.

9         Q.    And what you just described, was that

10   happening in connection with the pretraining of

11   Llama 3?

12        A.    That is my understanding.

13        Q.    Do you know if it was taking place before

14   Llama 3 was pretrained or after Llama 3 was

15   pretrained?

16             ATTORNEY HARTNETT:  Objection to the form.

17             THE WITNESS:  So this took place in --

18   this took place during the process of deciding on

19   the pretraining data for Llama 3, is my

20   understanding.

21   BY ATTORNEY STOLER:

22        Q.    Okay.

Page 64

```
 1   text data within LibGen?

 2        A.    No.  I don't know.

 3        Q.    All right.  Have you ever heard of a

 4   dataset referred to as "Anna's Archive"?

 5        A.    Yes.  I've heard of it.

 6        Q.    Is Anna's Archive a books dataset as well?

 7             ATTORNEY HARTNETT:   Object to the form.

 8             THE WITNESS:   So my understanding is that

 9   Anna's Archive is a superset of LibGen and contains

10   a number of things.  And so I -- my understanding is

11   that it contains books, but it -- it also has a lot

12   of other things in there.

13   BY ATTORNEY STOLER:

14        Q.    Do you know if Meta ever used the Books3

15   dataset to train any of its Llama models?

16        A.    Well, I've seen it mentioned in the Llama

17   paper; so for Llama 1.  I don't know what data was

18   trained for Llama 2.

19             So I -- there was a list of code names for

20   Llama 3, and so I -- I don't actually know which of

21   the code names corresponds to Llama 3.  So it may be

22   in there, but I don't know for certain.
```

12/11/2024          Richard Kadrey, et al. v. Meta Platforms, Inc.          Sean Bell
Highly Confidential - Attorneys' Eyes Only

Page 97

1    anyone in your team regarding the importance of

2    training on the Anna's Archive dataset?

3              ATTORNEY HARTNETT:  Objection to form.

4              THE WITNESS:  That's pretty vague.  What

5    do you mean by that?

6    BY ATTORNEY STOLER:

7         Q.    Has anyone on your team ever indicated to

8    you that it's important that Llama 4 train on the

9    Anna's Archive dataset?

10        A.    I mean, that would've been connected with

11   us expanding from LibGen to Anna's Archive.

12              And so the general understanding was that

13   we need more reasoning data.  "Reasoning" broadly

14   meaning STEM and mathematics and that one of the

15   best sources for reasoning data was in LibGen.  And

16   so for that reason, you know -- then there was a

17   question of, you know, will we use all of LibGen or

18   not and do we need to look at other sources or

19   related things similar to LibGen which was around

20   the conversation of saying, you know, let's also

21   think about Anna's Archive.

22        Q.    My question is just whether you remember

Page 131

 1   mean?

 2        A.    Like, downloading it from publicly

 3   available sources --

 4        Q.    Okay.

 5        A.    -- in the context of this bullet here.

 6        Q.    And can you think of an example of a

 7   publicly available dataset of books/PDFs that you

 8   downloaded?

 9        A.    Sure.  I mean, there is a whole bunch of

10   PDFs on Spidermate and web crawling.  And so, you

11   know, we made sure that when we're working with that

12   team that we ask them specifically to make sure to

13   prioritize downloading PDFs.

14        Q.    Okay.  And how is -- so I see that

15   sourcing books and PDFs is separate from licensing.

16            The -- is there -- do you know if there

17   was any attempt to license any of these books or PDF

18   datasets?

19        A.    So --

20            ATTORNEY HARTNETT:  Object to the form.

21            THE WITNESS:  So at this time, we -- we

22   had -- if I remember correctly, we had not really

Page 132

1    gotten into trying to license anything at scale.

2    And, you know, first we needed budget and sort of

3    leadership alignment around the very idea that --

4    should we be approaching publishers or not.

5              Since then, we started approaching some

6    publishers, and what we found was that the volume of

7    content that they have is orders of magnitude

8    smaller than what we needed for pretraining; and so

9    that changed the focus from "Okay.  Should we think

10   about licensing as a part of strategy of our

11   pretraining?"

12             And so there was a shift in thinking from

13   the time when we wrote this to now, which is that,

14   "Hey, licensing.  We should think of it as a --

15   something for posttraining."

16             Like, it just doesn't exist in the scale

17   that you need for pretraining to even be meaningful.

18   Like, whether or not you add -- like, say

19   hypothetically you licensed every book from some

20   publisher and tried to put it into pretraining.  It

21   would be a statistical noise that would be not even

22   measurable of whether or not you did it --

Page 133

1    pretraining scale, in terms of the volume of the

2    content.

3              And so -- you know, so since then, we've

4    been focusing on licensing for posttraining.  And so

5    because I'm focused mostly on pretraining, I have

6    not been involved -- that involved in licensing.

7    BY ATTORNEY STOLER:

8         Q.    When you say that it would be -- sorry.

9    When you say that it would amount to "statistical

10   noise" to include every book in the pretraining

11   data, what do you mean by that?

12        A.    So what I mean is that the minimum size

13   right now we're able to measure, approximately

14   speaking, is about 100 billion tokens.  So the final

15   data mix is, you know, say, 30 trillion in some of

16   the models and 60 trillion we're working towards for

17   the Llama 4 flagship model.

18              And the benchmarks themselves already are

19   a little bit noisy.  So there's already some plus or

20   minus on certainty.  And so if you change the data,

21   you have to know whether or not the change in data

22   resulted in a benchmark movement larger than the

Page 136

```
 1       A.    Well, I mean, the scale that we're able to

 2    get in terms of PDFs -- and there's a lot of books

 3    in those PDFs -- is something that we can measure,

 4    you know, in terms of what you can download in terms

 5    of publicly available sources.

 6               So I'm talking about specifically a

 7    licensing strategy where there is a bit of a song

 8    and dance with, you know -- so say we go to someone

 9    who has content to license and, you know, there's

10    a -- you know, haggling process on the price.  You

11    know, it takes a bunch of their time; it takes our

12    time.  They might send us a sample of the data.

13    They don't necessarily want to the give us all data

14    because then they think, "Well, I've just given it

15    to you."

16               So now -- so the problem is compounded

17    even more that they might give us a tiny sample of

18    data, which is very hard to figure out what it -- so

19    the problem of figuring out the value is made 100

20    times harder because they give us a tiny sample.

21               And as I mentioned earlier, what we know

22    is that human intuitions studying an individual
```

Page 304

1                CERTIFICATE OF SHORTHAND REPORTER

2

3        I, Michael P. Hensley, Registered Diplomate

4    Reporter for the State of California, CSR No. 14114,

5    the officer before whom the foregoing deposition was

6    taken, do hereby certify that the foregoing

7    transcript is a true and correct record of the

8    testimony given; that said testimony was taken by me

9    stenographically and thereafter reduced to

10   typewriting under my direction; that reading and

11   signing was not requested; and that I am neither

12   counsel for, related to, nor employed by any of the

13   parties to this case and have no interest, financial

14   or otherwise, in its outcome.

15

16

17

18

19        Michael P. Hensley, CSR, RDR

20

21

22